# Advancements in the human genome reference assembly (GRCh38)

Tayebeh Rezaie, Ph.D.
NCBI
16 October 2019

# Genome Reference Consortium

**GRC**

- Valerie Schneider
- Kerstin Howe
- Tina Graves
- Paul Flicek
- Tayebeh Rezaie
- Nathan Bouk
- Hsiu-Chuan Chen
- Jo Wood
- Joanna Collins
- Sarah Pelan
- Will Chow
- James Torrance
- Ying Sims
- Derek Albracht
- Milinn Kremitzki

**Thanks to many GRC Collaborators**
https://www.ncbi.nlm.nih.gov/grc/credits/

# History of reference assembly

**GRCh38/Reference genome:**

- A critical resource to the basic & clinical research community, coordinate system, annotation source & discovery of disease-associated variants
- Sanger seq. clone-based from **H**uman **G**enome **P**roject; multiple individuals
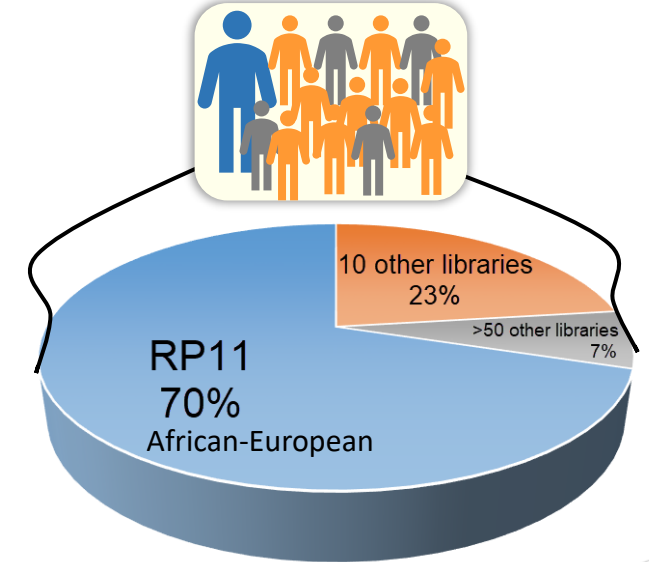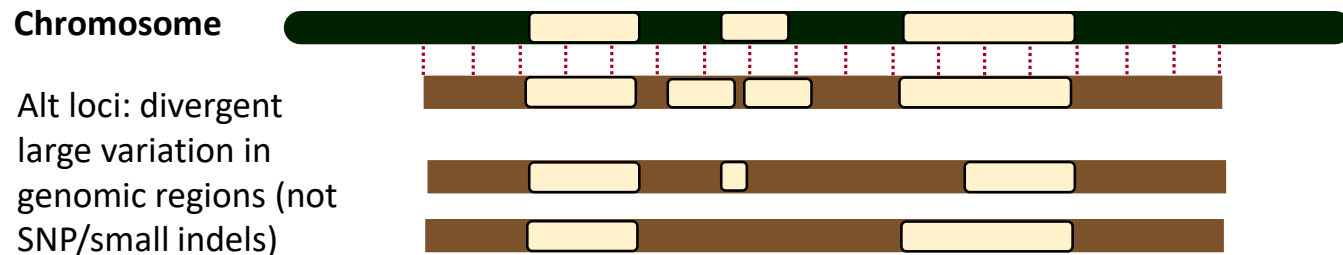
**Mosaic haploid**  Individual 1   Individual 2   Individual 1

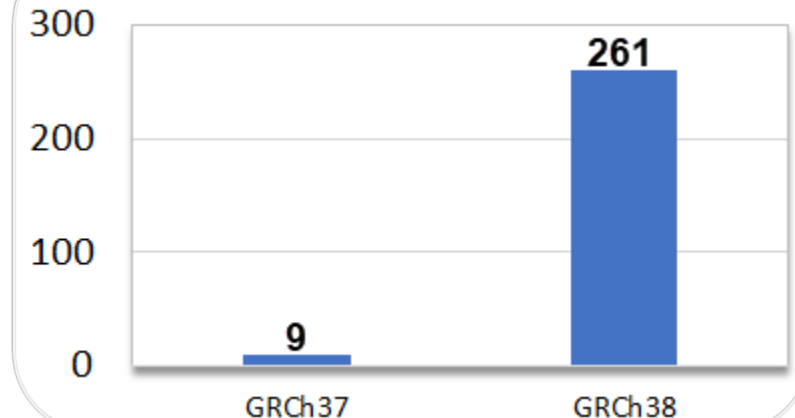**HGP** ➡ **GRC: reference maintaining, improving and updates**

**HGP model (2003):** each genomic region was represented with one sequence

Chromosome

**Current model:** ALT LOCI added to represent population genomic diversity

Chromosome

Alt loci: divergent large variation in genomic regions (not SNP/small indels)

10 other libraries
23%

RP11
70%
African-European

>50 other libraries
7%

## Number of ALT LOCI



300

261

200

100

9

0

GRCh37    GRCh38

# Reference assembly updates



Resolved Since GRCh38 n=430

- Variation 46.5%
- Seq. error 21%
- Gap 15%
- GRC 10.5%
- Path 3.5%
- Missing seq. 1.6%
- Unknown 0.5%
- Localization 1.4%

**Reference updates released as patches: 185/430 (42%)**

**The new version of the reference should capture ALL the updates to GRCh38**

Legend:
◀ Alt Loci
● Fix Patch
● Novel Patch

- **Major/coordinate-changing: GRCh38 (Dec 2013)**
- **Patches/no coordinate-change: GRCh38.p13 (Mar 2019)**

- 113 Fix patches: Add >3.88 Mb novel seq
- 72 Novel patches: Add >1.1 Mb novel seq
- 261 ALT Loci: Add 3.6 Mb novel seq

**The notion for variant representation has started long time ago.**
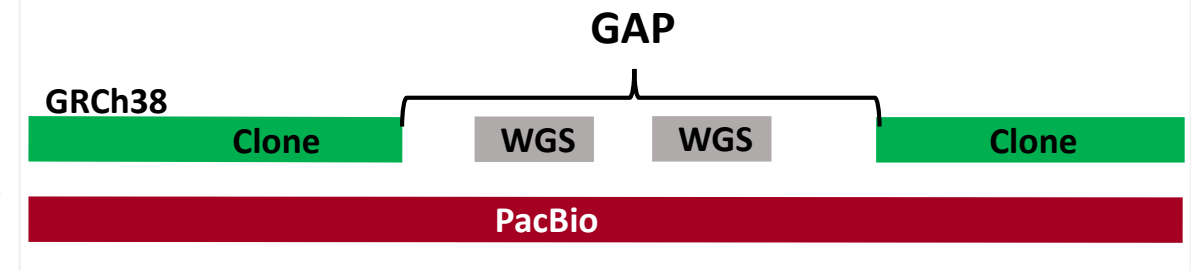
# Curation of reference assembly

- Issue sources: GRC assembly evaluation, reports from collaborators, community, literature
- Technology: sequencing, FISH, Optical Mapping
- Data resources: sequences generated by GRC or available in public database (clones, WGS, PCR products)

## Evaluation of gaps in GRCh38

- Gap count = 196

Excluded biological gaps & gaps within WGS scaffolds
- Reports of new assm that can close ref. gaps
- To identify gaps that can be spanned



**Alignments of 8 diploid PacBio assemblies to the reference:**
- Spanned with the same amount of seq: 26 (missing seq.)
- Spanned with varying amount of seq: 3 (variation)
- Spanned by some not all assemblies: 24 (complex, missing + variation)
- The remaining gaps are under review

https://www.genome.wustl.edu/research/projects/

# Curation of reference assembly: Missing sequences
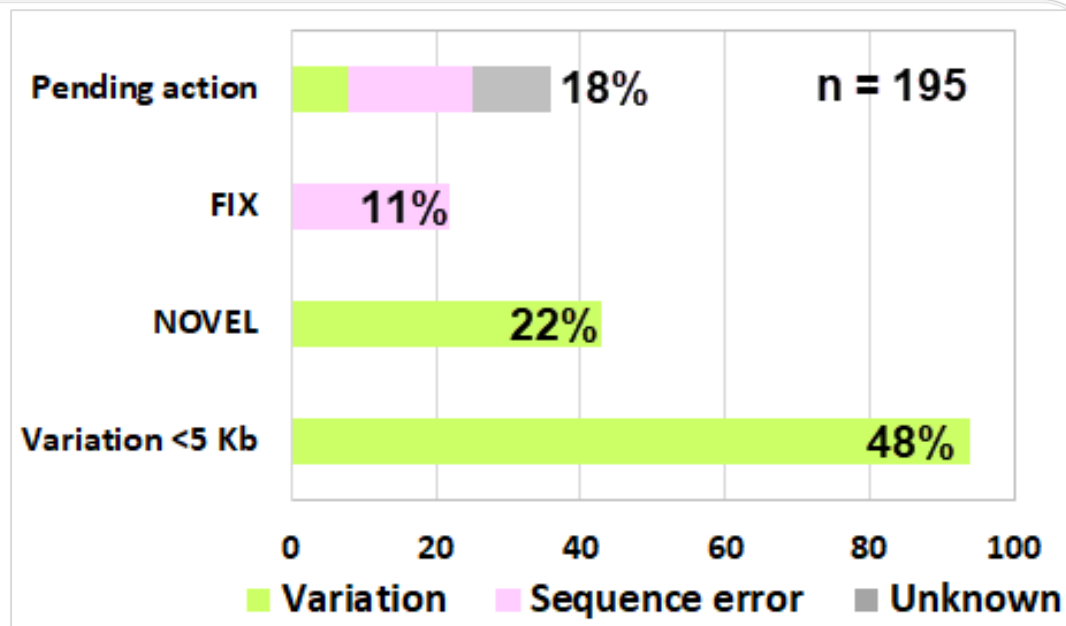
Evaluation to distinguish error vs. variation → Find chr. context for missing seq.

→ Add variants (>5 kb) as novel patches

Data sources:
- Eichler's lab (Kidd et al. (2010) PMID: 20440878), structurally variant fosmid seq.
- Heng Li (GCA_000786075.2), a set of non-redundant seq. absent in GRCh38 and ALTs

**Reported genome issues = 195**
- Resolved no change: 94 (variation < 5 Kb)
- Patches (started adding from p1 in 2014)
  - FIX = 22
  - NOVEL = 43
- Pending action: 36 (Variation 8, sequence error 17, Unknown 11)



| | |
|---|---|
| Pending action | 18% n = 195 |
| FIX | 11% |
| NOVEL | 22% |
| Variation <5 Kb | 48% |

Variation ■ Sequence error ■ Unknown

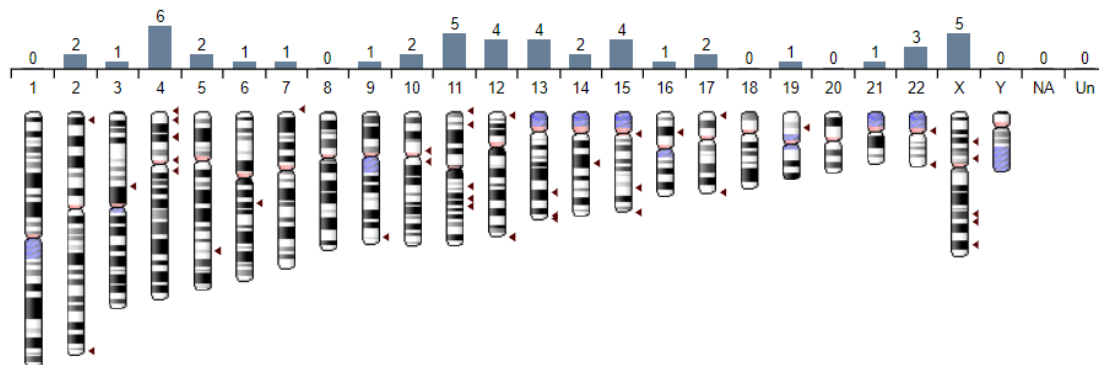# GRCh38.p13 updates to reference assembly

**The most recent curation to GRCh38:**

- FIX patches (43) + NOVEL (2)
- Added >0.5 Mb novel sequence
    - Gap closure: 28
    - Seq. error correction: 8
    - Path: 2
    - For p-arm of acrocentric chrs: 5

- **Highlights of p13:**
    - Improved clinically important genomic regions
        - Prader-Willi (5.5 Mb, 1.63 Mb unique)
        - CT47A gene cluster
    - Improved gene representations: SLC5A11, GCNT2, SAMD1, GRCK1, C1R, ECSCR, 5S rRNA

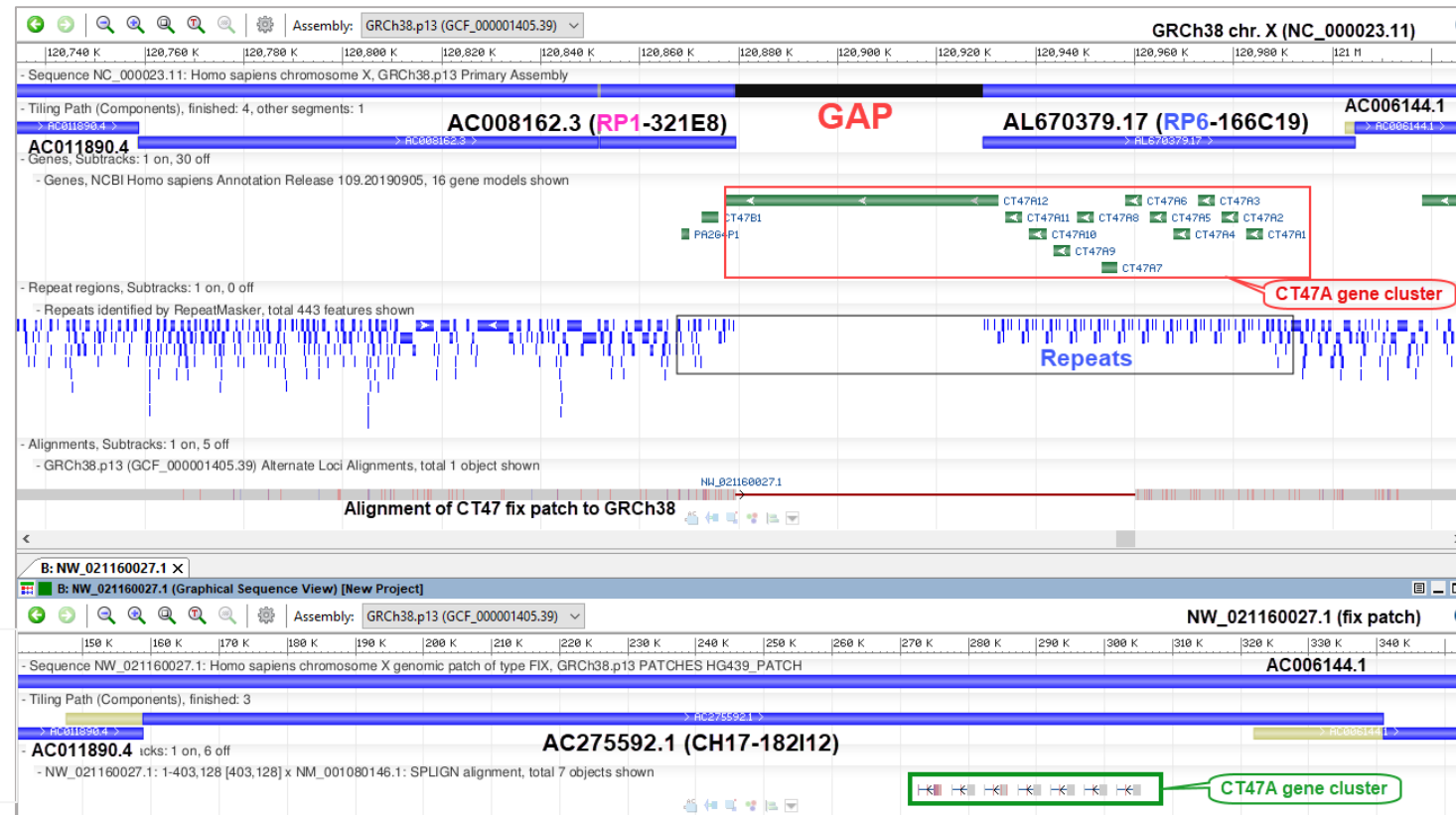**Chr. distribution of GRCh38.p13 patches**



- Sequence data sources for updates:
    - CHM1 assm: 21
    - CHM13 assm: 12
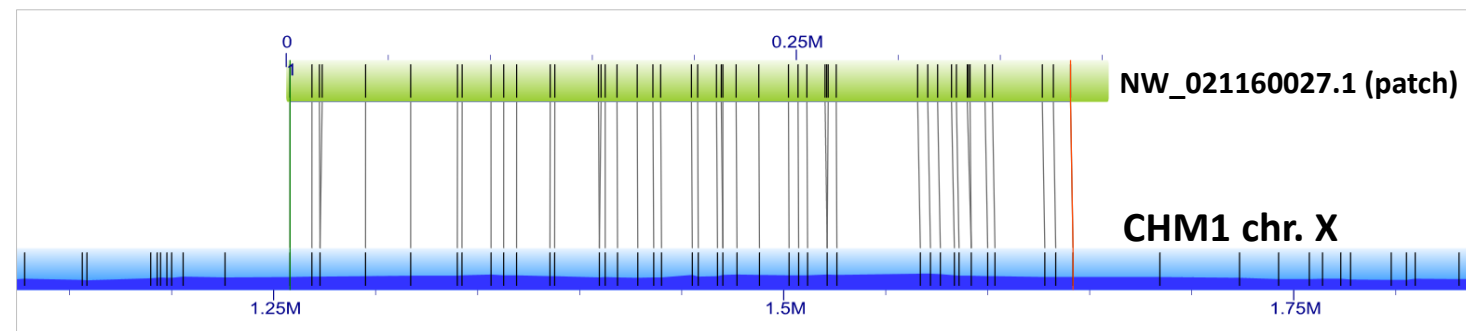    - Other WGS assm: 3
    - Clones: 9

# Correction of an assembly false gap caused by haplotype incompatibility

**Mix haplotype representation of CT47A in GRCh38**
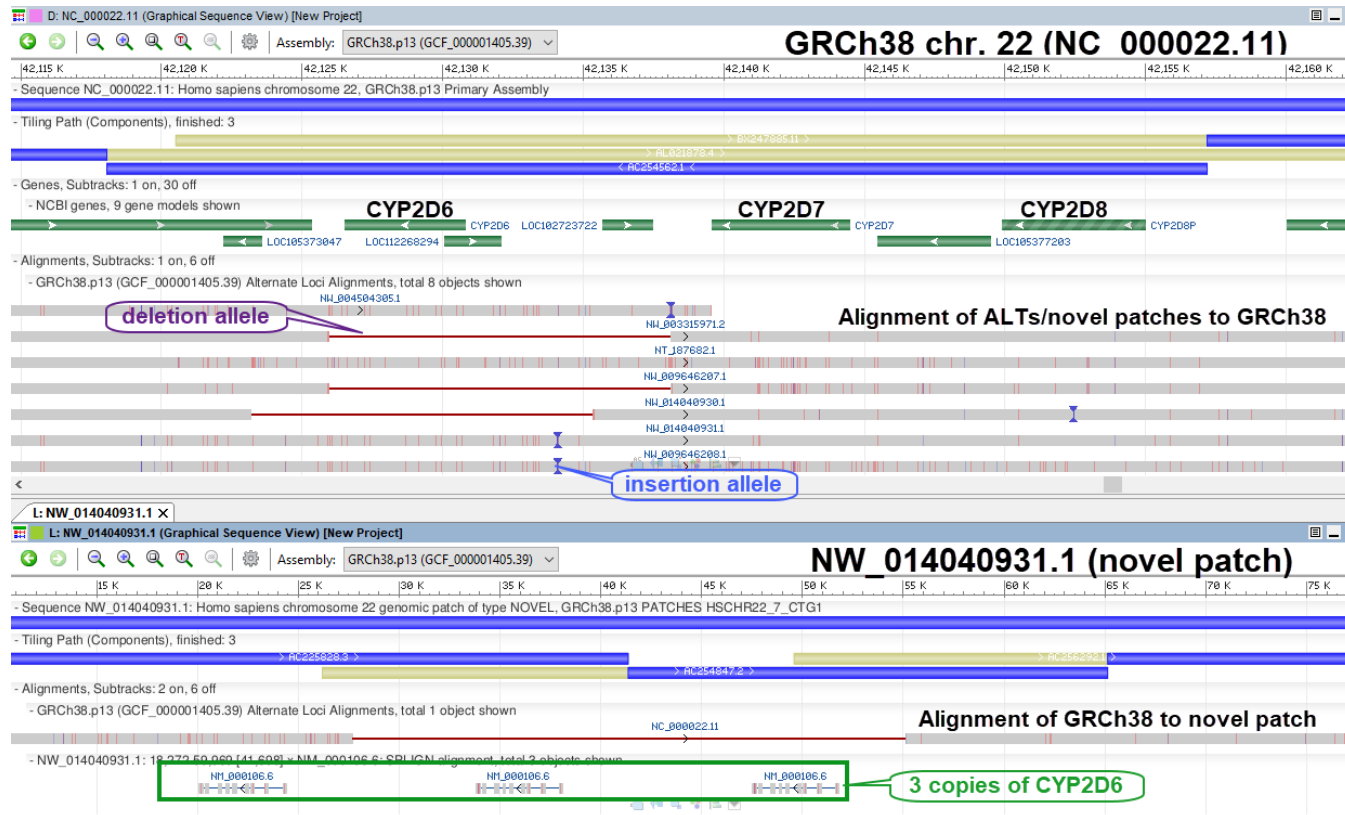**Long haplotype: 12 copies**

**Single haplotype representation of CT47A in GRCh38.p13**
**Short haplotype: 7 copies**

**CHM1 Optical Map supporting the updated CT47A haplotype**

# CYP2D6 haplotypes: genomic diversity of a clinically important region
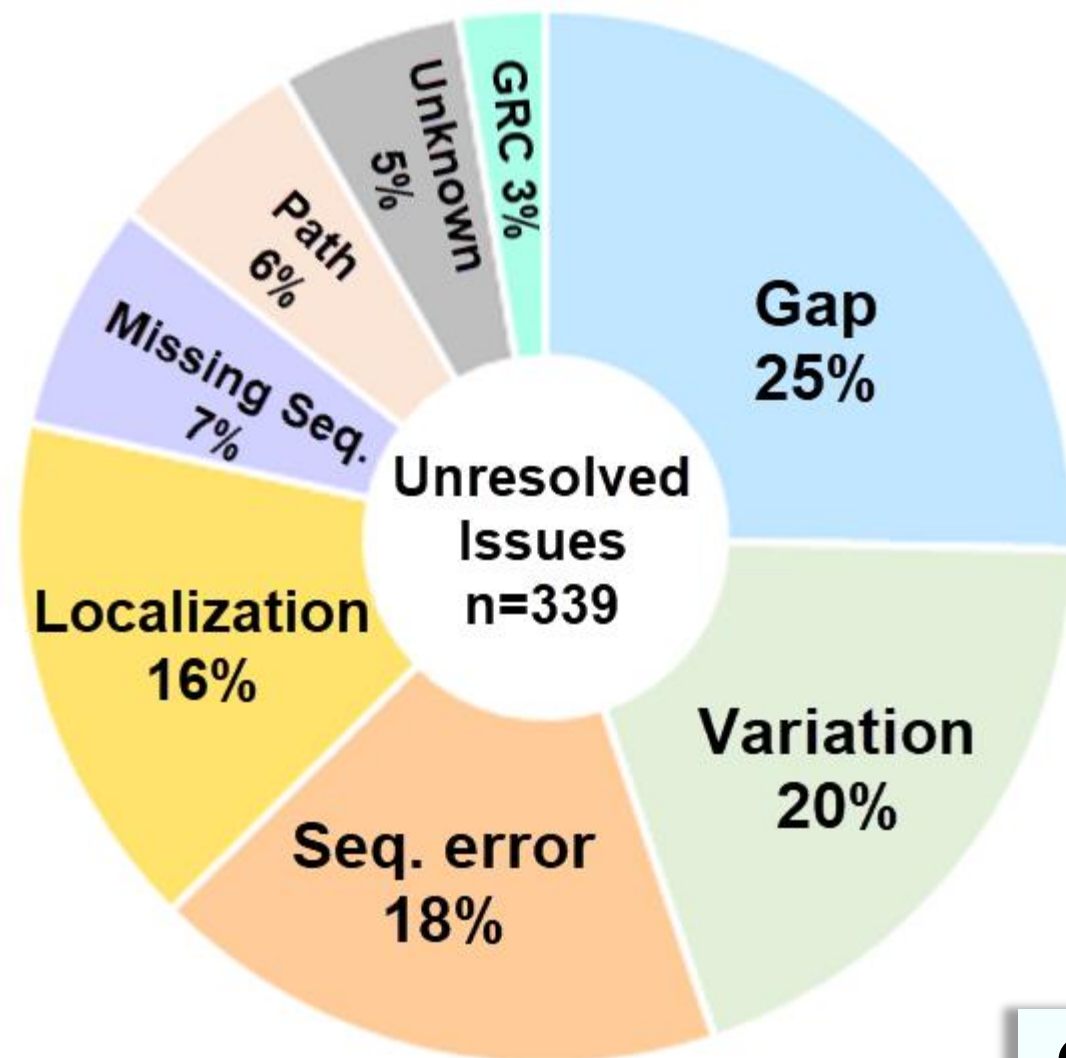# Involved in metabolizing many prescribed drugs

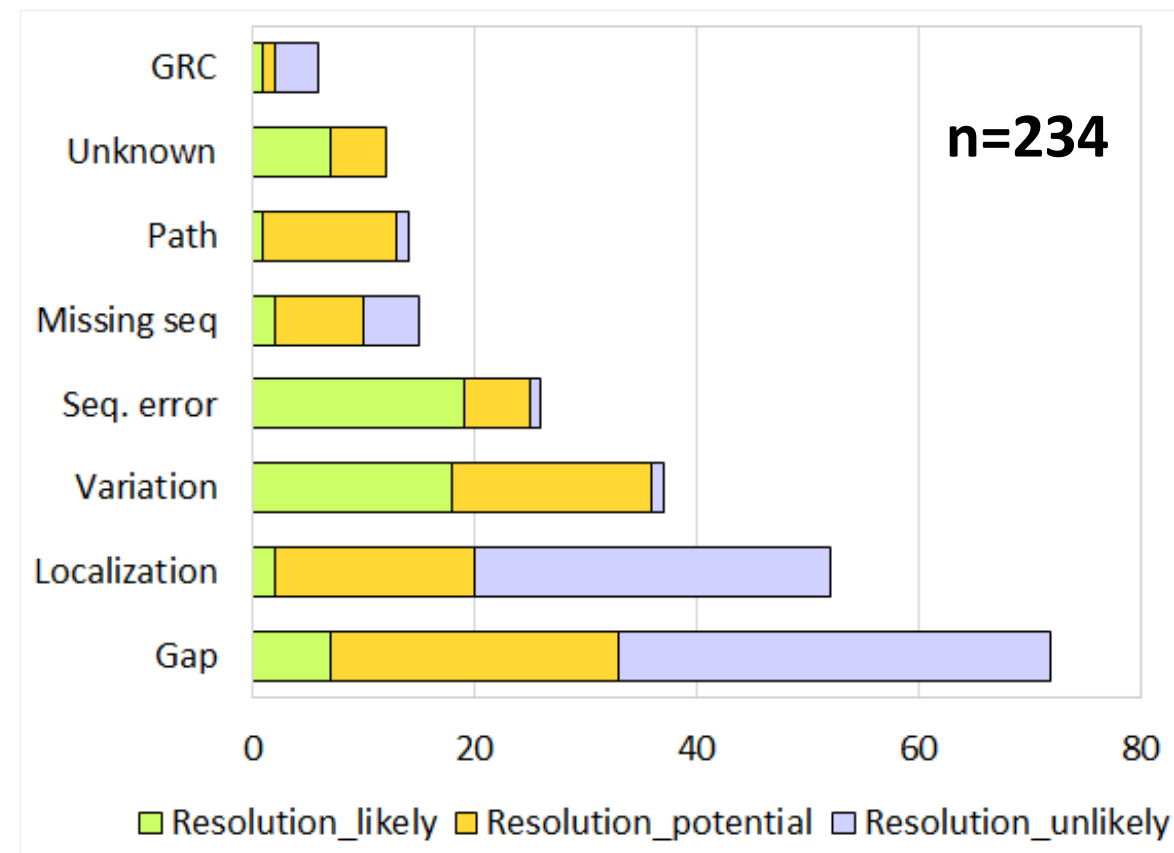**Alignment of alt loci and patch scaffolds to the CYP2D6 region of chr. 22**

**Scaffolds providing alternate sequence representations of CYP2D6 region**



| GenBank Acc. | RefSeq Acc. | Population | CYP2D6 | CYP2D7 | CYP2D8 |
|---|---|---|---|---|---|
| KN196485.1 | NW_009646207.1 | African | Deletion | Single Copy | Single Copy |
| KB663609.1 | NW_004504305.1 | African | Duplication | Single Copy | Single Copy |
| KN196486.1 | NW_009646208.1 | East Asian | Duplication | Single Copy | Single Copy |
| KQ458387.1 | NW_014040930.1 | East Asian | Deletion | Single Copy | Duplication |
| KQ458388.1 | NW_014040931.1 | East Asian | 3 Copies | Single Copy | Single Copy |
| KQ759761.1 | NW_015148968.1 | European | Single Copy | Duplication | Single Copy |
| GL383582.2 | NW_003315971.2 | Unknown | Deletion | Single Copy | Single Copy |

# Genome Reference Consortium

## Unresolved genome issues



Unresolved Issues n=339

- Gap 25%
- Variation 20%
- Seq. error 18%
- Localization 16%
- Missing Seq. 7%
- Path 6%
- Unknown 5%
- GRC 3%

## Current curation status



n=234

Resolution likelihoods as determined by the GRC review

□ Resolution_likely  □ Resolution_potential  □ Resolution_unlikely

## GRCh38.p14 is planned for release in 2020

# Conclusion and Future

*The future is* **BRIGHT**

## GRCh38.p14: coming in 2020

## The reference has informed its own evolution.



Human Pangenome Reference Consortium

HPSC  PRR  PRT  RST  SAB  NHGRI  Steering Committee

Human Pangenome Reference Center

Outreach and Education

WashU

UCSC — EBI

HPRP Logistical Coordination Center

Executive Committee

Assembly Validation

Variant Calling

Pan-genome Framework

Pan-genome Maintenance

Global Scientific & Clinical Community

Consortium Partners
GA4GH, GIAB, ClinVar, HGSVC, ClinGen, ENCODE, 4DN, AnVIL, CCDG, CMG, TOPMed, *All of Us*

Genome Reference Consortium

MGI, a GRC member, has been awarded by NHGRI to:
- Produce 350 whole genome phased diploid assm.;
- Identify SVs between samples and current GRCh38;
- Incorporate those SVs into the Reference, likely as a graph representation.

GRCh39 is pending. The GRC is engaged in validation, providing curation tools and support to the pan-genome assemblies.